

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 11-184825

(43)Date of publication of application : 09.07.1999

(51)Int.Cl. G06F 15/16
G06F 11/20
G06F 13/00

(21)Application number : 09-350391

(71)Applicant : MITSUBISHI ELECTRIC CORP

(22)Date of filing : 19.12.1997

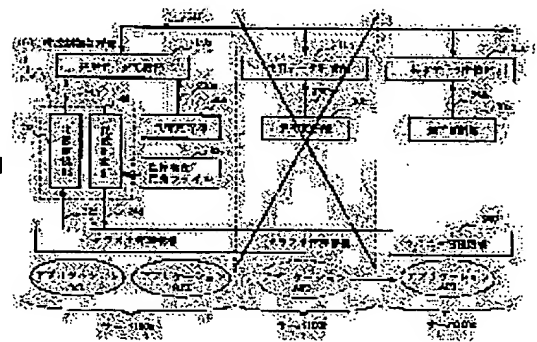
(72)Inventor : YAMANISHI HIROYUKI
SAITO AKIO

(54) CLUSTER SYSTEM

(57)Abstract:

PROBLEM TO BE SOLVED: To dynamically start/stop an application, decide a start computer, decide a succeeding calculator in accordance with the state change of the computer and the application and to optimize a cluster system in the cluster system.

SOLUTION: A state monitoring part 20 and a load update part 30 take in information showing the state of the computer and the application. Information which is taken in is managed by a common data storage part 10. A constitution control part 40 outputs the instruction of the succession of the application and the change of the succeeding computer to a cluster management mechanism 900 when it receives a trigger from a state monitoring part 20 and a load update part 30.



LEGAL STATUS

[Date of request for examination] 02.02.1998

[Date of sending the examiner's decision of rejection] 22.05.2001

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

Copyright (C); 1998,2003 Japan Patent Office

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平11-184825

(43) 公開日 平成11年(1999) 7月9日

(51) Int.Cl.⁶

識別記号

F I

G 0 6 F 15/16

3 8 0

G 0 6 F 15/16

3 8 0 Z

11/20

3 1 0

11/20

3 1 0 F

13/00

3 5 5

13/00

3 5 5

審査請求 有 請求項の数 8 O L (全 19 頁)

(21) 出願番号

特願平9-350391

(22) 出願日

平成9年(1997)12月19日

(71) 出願人 000006013

三菱電機株式会社

東京都千代田区丸の内二丁目2番3号

(72) 発明者 山西 宏幸

東京都千代田区丸の内二丁目2番3号 三
菱電機株式会社内

(72) 発明者 斎藤 彰男

東京都千代田区丸の内二丁目2番3号 三
菱電機株式会社内

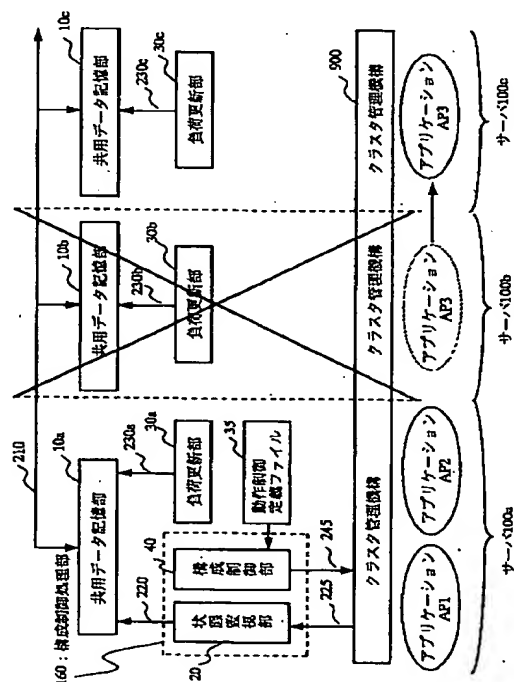
(74) 代理人 弁理士 宮田 金雄 (外2名)

(54) 【発明の名称】 クラスタシステム

(57) 【要約】

【課題】 クラスタシステムにおいて、計算機及びアプリケーションの状態変化に応じてアプリケーションの起動、停止、起動計算機の決定及び引き継ぎ計算機の決定を動的に行い、クラスタシステムの最適化を図る。

【解決手段】 状態監視部20、負荷更新部30が計算機及びアプリケーションの状態を示す情報を取り込む。取り込んだ情報は、共用データ記憶部10で管理する。構成制御部40は、トリガを状態監視部20、負荷更新部30から受け取ると、アプリケーションの引き継ぎや引き継ぎ計算機の変更などの指示をクラスタ管理機構900に出す。



1

【特許請求の範囲】

【請求項1】 1つ以上のアプリケーションを実行可能な複数台のサーバをメンバーとして構成されるクラスタシステムであって、上記複数台のサーバのいずれか1台に上記アプリケーションを実行させるとともに、上記1台のサーバが上記アプリケーションを実行不可能である時に他のサーバに上記アプリケーションの実行を引き継がせるクラスタシステムにおいて、

以下の要素を有するクラスタシステム

(a) 上記複数台のサーバから共通して利用可能な共用データとして、上記複数台のサーバの状態を示すサーバ状態情報と上記複数台のサーバの負荷を示す負荷情報と上記アプリケーションの動作に関するアプリケーション動作情報とを記憶する共用データ記憶部、(b) 上記複数台のサーバの状態を取得して上記サーバ状態情報を更新する状態監視部、(c) 上記複数台のサーバの負荷を取得して上記負荷情報を更新する負荷更新部、(d) 上記共用データ記憶部を参照して、上記アプリケーションと上記複数台のサーバとの対応を動的に制御する構成制御部。

【請求項2】 上記共用データ記憶部は、上記クラスタシステムの制御に関する規約を定義する規約情報を記憶するとともに、上記構成制御部は、上記クラスタシステムの起動の際のクラスタ構成時に上記規約情報と上記サーバ状態情報とを参照して上記アプリケーション動作情報を作成することを特徴とする請求項1に記載のクラスタシステム。

【請求項3】 上記構成制御部は、上記クラスタシステムのクラスタ構成後に新たなサーバが起動してメンバーとして上記クラスタシステムに参加した時、上記サーバ状態情報と上記アプリケーション動作情報とを参照して上記アプリケーションの引継ぎを行うことを特徴とする請求項1に記載のクラスタシステム。

【請求項4】 上記構成制御部は、上記クラスタシステムのクラスタ構成後にメンバーの移動が発生した時、上記サーバ状態情報と上記アプリケーション動作情報とを参照して上記アプリケーションを引き継ぐ引継ぎサーバを決定して上記アプリケーション動作情報を更新することを特徴とする請求項1に記載のクラスタシステム。

【請求項5】 上記構成制御部は、上記複数のサーバのいずれかの負荷状態に変化が起こった時、上記サーバ負荷情報と上記アプリケーション動作情報とを参照して上記アプリケーションの引継ぎを行うことを特徴とする請求項1に記載のクラスタシステム。

【請求項6】 上記規約情報は、ユーザが設定する定義ファイルをもとに作成されることを特徴とする請求項2に記載のクラスタシステム。

【請求項7】 上記クラスタシステムは、パターン化されたテンプレートを予め用意し、用意したテンプレートのいずれかをもとに上記規約情報を作成することを特徴

2

とする請求項2に記載のクラスタシステム。

【請求項8】 上記負荷更新部は、所定のタイミングで定期的により上記複数台のサーバの負荷を取得して上記負荷情報を更新し、上記構成制御部は、上記負荷情報を参照してアプリケーション引継ぎに際し、上記アプリケーションを引き継ぐ引継ぎサーバを決定して上記アプリケーションの引継ぎを行うことを特徴とする請求項1に記載のクラスタシステム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】この発明は、分散した複数の計算機を一まとまりにして構成したクラスタシステムに関するものである。特に、クラスタシステムでのアプリケーションの引き継ぎに関するものである。

【0002】

【従来の技術】クラスタシステムにおいては、複数の計算機で1つ以上のアプリケーションを実行し、アプリケーション、もしくはアプリケーションを実行中の計算機が異常停止した場合、クラスタ内の他の計算機でアプリケーションを引き継いで実行する。この、アプリケーションを引き継いで実行する計算機を引継ぎ計算機と呼ぶ。引継ぎ計算機は、クラスタシステム内で、予め設定される。このような、アプリケーションの切り替えを示した従来の技術として、特開平6-28206号公報に開示された「クラスタシステムのデータ処理ステーション障害時の復旧方式」がある。この技術は、運用側の故障により無条件に予備側に切り替えるものである。だが、その後の新しい運用側の障害には対処できないという欠点があった。

【0003】

【発明が解決しようとする課題】このように、従来のクラスタの構成制御方式では、自動引継ぎ制御しか行われておらず、異常停止した計算機が復帰した等の事象によってクラスタ構成計算機メンバーが変更された場合、アプリケーションの引継ぎ計算機の再設定などは、システム管理者がコマンド等を実行して行う必要があった。また、アプリケーションが計算機に引き継がれる際、計算機やアプリケーションの状況に対応した引継ぎ計算機の自動変更処理が行われておらず、重複障害発生時に1つの計算機にアプリケーションが集中してしまうことなどもある。

【0004】この発明は、上記のような問題点を解決するためになされたものであり、計算機及びアプリケーションの状態変化に応じて、アプリケーションの起動、停止、起動計算機の決定を自動的に行うクラスタシステムを得ることを目的としている。また、計算機及びアプリケーションの状態変化に応じて、動的に引継ぎ計算機の決定を行い、クラスタ構成を最適に制御するクラスタシステムを得ることを目的としている。また、予め形式の決まった運転形態を簡単に設定することができる機能を

10

20

30

40

50

3

備えたクラスタシステムを得ることを目的としている。

【0005】

【課題を解決するための手段】この発明に係るクラスタシステムは、1つ以上のアプリケーションを実行可能な複数台のサーバをメンバーとして構成されるクラスタシステムであって、上記複数台のサーバのいずれか1台に上記アプリケーションを実行させるとともに、上記1台のサーバが上記アプリケーションを実行不可能である時に他のサーバに上記アプリケーションの実行を引き継がせるクラスタシステムにおいて、以下の要素を有することを特徴とする。

(a) 上記複数台のサーバから共通して利用可能な共用データとして、上記複数台のサーバの状態を示すサーバ状態情報と上記複数台のサーバの負荷を示す負荷情報と上記アプリケーションの動作に関するアプリケーション動作情報とを記憶する共用データ記憶部、(b) 上記複数台のサーバの状態を取得して上記サーバ状態情報を更新する状態監視部、(c) 上記複数台のサーバの負荷を取得して上記負荷情報を更新する負荷更新部、(d) 上記共用データ記憶部を参照して、上記アプリケーションと上記複数台のサーバとの対応を動的に制御する構成制御部。

【0006】上記共用データ記憶部は、上記クラスタシステムの制御に関する規約を定義する規約情報を記憶するとともに、上記構成制御部は、上記クラスタシステムの起動の際のクラスタ構成時に上記規約情報と上記サーバ状態情報とを参照して上記アプリケーション動作情報を作成することを特徴とする。

【0007】上記構成制御部は、上記クラスタシステムのクラスタ構成後に新たなサーバが起動してメンバーとして上記クラスタシステムに参加した時、上記サーバ状態情報と上記アプリケーション動作情報とを参照して上記アプリケーションの引き継ぎを行うことを特徴とする。

【0008】上記構成制御部は、上記クラスタシステムのクラスタ構成後にメンバーの移動が発生した時、上記サーバ状態情報と上記アプリケーション動作情報とを参照して上記アプリケーションを引き継ぐ引継ぎサーバを決定して上記アプリケーション動作情報を更新することを特徴とする。

【0009】上記構成制御部は、上記複数のサーバのいずれかの負荷状態に変化が起こった時、上記サーバ負荷情報と上記アプリケーション動作情報とを参照して上記アプリケーションの引き継ぎを行うことを特徴とする。

【0010】上記規約情報は、ユーザが設定する定義ファイルをもとに作成されることを特徴とする。

【0011】上記クラスタシステムは、パターン化されたテンプレートを予め用意し、用意したテンプレートのいずれかをもとに上記規約情報を作成することを特徴とする。

【0012】上記負荷更新部は、所定のタイミングで定

4

期的に上記複数台のサーバの負荷を取得して上記負荷情報を更新し、上記構成制御部は、上記負荷情報を参照してアプリケーション引き継ぎに際し、上記アプリケーションを引き継ぐ引継ぎサーバを決定して上記アプリケーションの引き継ぎを行うことを特徴とする。

【0013】

【発明の実施の形態】実施の形態1. 最初に、この発明のクラスタシステムが実現する機能の概要を説明する。

1. 実現機能の概要

(a) 計算機の初期起動時、即ち、個々の計算機が電源投入され、クラスタの初期構成が行われる時、クラスタ内のアプリケーションの引き継ぎ計算機を自動制御し、決定する機能。

(b) 計算機の起動時、アプリケーションを自動的に引き継がせる機能。

(c) 計算機異常停止時又は異常停止からの復帰時など、クラスタを構成する計算機に変更があった際、クラスタ内のアプリケーションの引き継ぎ計算機を自動制御して変更する機能。

(d) 計算機が通常、起動している状態(CPU、メモリ、ネットワーク)の負荷を監視し、その状態の変化に応じて、アプリケーションを自動制御する機能。

(e) プログラミングを行うことなく、定義ファイルを設定するだけで、構成の再構築制御を行う機能。

(f) ホットスタンバイ、ロードシェア、N:1バックアップ型など、予め形式の決まった運転形態を簡単に設定することができる機能。

(g) アプリケーションを自動で引き継がせる場合、各計算機の負荷も考慮し、引き継ぎ計算機を決定する機能。

【0014】次に、図を用いてこの発明のクラスタシステムの構成を説明する。図1は、この発明のクラスタシステムの構成の一例を示すブロック図である。図において、900は既存のクラスタ管理機構である。この発明のクラスタシステムでは、サーバ100a~100c(計算機ともいう)の状態及びアプリケーションの状態に関する情報をサーバ100aの状態監視部20でクラスタ管理機構900から取り込み、それらの情報を共用データ記憶部10aに記憶する(矢印220)。この共用データ記憶部10aは、矢印210に示すように、他のサーバの共用データ記憶部10b、共用データ記憶部10cにコピーして同じ状態に置くよう管理されるものとする。また、負荷更新部30a~30cは、サーバ100a~100cの負荷に関する情報を得て共用データ記憶部10aに記憶する(230a~230c)。構成制御部40は、システムの状態変化を状態監視部20から受けると、又はシステムの負荷の変化を各計算機に存在する負荷更新部30a~30cから受けると、アプリケーションの引き継ぎや、引き継ぎ計算機の変更などの指示をクラスタ管理機構900へ出す。状態監視部20と構成制御部40を構成制御処理部160と呼ぶ。ここで

は、構成制御処理部160がサーバ100a上で動作する場合を示しているが、サーバ100bやサーバ100c上で動作しても構わない。

【0015】以下に、前述したクラスタシステムの機能を実現する実現手段について説明する。

2. 実現手段の概要

(a) クラスタの初期起動の通知を受けた構成制御部40は、実行中のアプリケーションの引継ぎ計算機を、アプリケーションを実行していない計算機に設定する。

(b) ある計算機が動作中のクラスタに参加してきた時、その通知を受けた構成制御部40は、その計算機上で動作すべきアプリケーションを検索し、動作すべきアプリケーションが検索されれば、そのアプリケーションを引き継がせる指示を出す。

(c) クラスタに参加中の計算機が停止したり、動作中のクラスタにある計算機が参加した場合、その通知を受けた構成制御部40は、引継ぎ計算機を変更する指示を出す。

(d) 各計算機の負荷状態(CPU、メモリ、ネットワーク)を監視し、その状態を管理する機能を構成制御部40に持たせる。実行中のアプリケーションをその負荷の状況に応じて、引継ぎ計算機に引き継がせる指示を出す。

(e) 計算機が異常停止した時、停止していた計算機が起動した時のアプリケーションの動作について定義した制御動作定義ファイル35を管理コマンドを用いて、各計算機へ配布/配置する。構成制御部40は、各要素に変化があった場合、その情報に従った自動制御を行う。

(f) 予めパターン化されたクラスタの自動構成制御を行うための定義情報のテンプレートを用意し、ユーザに開放する。テンプレートをもとに、制御動作定義ファイル35が作成される。

(g) 負荷更新部30a~30cが定期的に負荷の情報を取り込み、取り込んだ情報を共用データ記憶部10a~10cへと反映する。各計算機の負荷情報を管理する管理機能を持った構成制御部40が、その負荷情報が示す負荷状況に応じて最適なアプリケーションの引継ぎを指示する。図2に、アプリケーションの引継ぎの例を示す。図2は、サーバ100bが異常停止し、アプリケーションAP3がサーバ100cに引き継がれた状態を示している。

【0016】3. ハードウェア構成

図3は、この発明のクラスタシステムのハードウェア構成の一例を示す図である。ハードウェアは、各ローカルディスク200a~200nを持つN台のサーバ100a~100n、各サーバ間の通信路であるネットワーク1000から構成される。

【0017】4. ソフトウェア構成

図4、図5は、マルチサーバ上のソフトウェアの位置付けを説明する図である。ソフトウェアは、大別して下記

部位から構成される。

(a) 既存のクラスタ管理機構900

アプリケーションパッケージを、定義された情報に従ってサーバに割り付け、実行させるクラスタ管理機構である。クラスタ管理機構900は、アプリケーション実行中のサーバが異常停止した場合、アプリケーション毎に定義されている別サーバで再起動するなどの制御を行う。本実施の形態では、ヒューレット・パッカード社(米国のHewlett-Packard Company、HP社ともいう。Hewlett-Packard Company、HPはHP社の商標又は登録商標)のMC/ServiceGuardをクラスタ管理機構として使用する場合を例にとって説明する。クラスタ管理機構は、既存の製品であり、本発明以前の従来からある技術である。

【0018】(b) 構成制御処理部160

クラスタを構成する全てのサーバ上で動作可能な構成制御処理部160がある。構成制御処理部160は、状態監視部20と構成制御部40からなる。構成制御処理部160は、クラスタを構成する全てのサーバのうち、いずれかのサーバ上で動作する。構成制御処理部160は、全てのサーバで動作可能なもので、クラスタ管理機構900に対して予め設定することにより、下記のようにクラスタ管理機構900によって起動される。

・初期起動

クラスタ管理機構900の設定(起動サーバの優先順位指定)に従い、現状構成内の高優先順位のサーバで起動される。

・動作サーバ停止時の再起動

動作中のサーバが異常停止した場合、クラスタ管理機構の設定に記述された次優先順位のサーバで再起動され、引き続き動作する。図5は、サーバ100aが停止し、構成制御処理部160がサーバ100bで引き続き動作する場合を示している(矢印219)。この初期/再起動のメカニズムは、既存のクラスタ管理機構が行う。構成制御処理部160は、前述した機能を実現し、アプリケーションの移動先サーバの決定などを行う。決定された移動先のサーバの指示は、クラスタ管理機構900に対して行われる。また、アプリケーションの起動、移動等の制御は、クラスタ管理機構900により行われる。クラスタ管理機構900は、構成制御処理部160からの指示に従い動作する。即ち、この実施の形態の構成制御処理部160は、クラスタ管理機構900の存在を前提としている。構成制御処理部160は、直接アプリケーションを起動したり、アプリケーションの引き継ぎを行ったりすることはない。

【0019】(c) 共用データ記憶部10a~10c

ネットワークを介した共用メモリ機能(あるサーバで書き込んだ内容を別サーバに反映することで、別サーバがその内容を参照できる機能)をサポートする記憶部であ

る。構成制御処理部160が別サーバで再起動された場合の引き継ぎデータとして使用される。この共用メモリ機能の実現方式は問わないものとする。この実施の形態では、この共用メモリ機能にサポートされた共用データ記憶部を用いる場合について説明する。データの共用方式には、特殊な装置や共用ディスクを使用した様々な方式が考えられるが、本実施の形態は、ネットワークを介した方式で記述する。

【0020】(d) アプリケーションプロセス(アプリケーションAP1~AP3)

構成制御処理部160の動作結果として、障害時のバックアップサーバ設定などが行われるユーザのアプリケーションプログラム。実際のサーバ起動は、クラスタ管理機構が行う。本発明のクラスタシステムは、クラスタ管理機構900に対して指示を出す。図5の矢印234は、アプリケーションAP1、アプリケーションAP2がサーバ100aの停止に伴い、サーバ100cで引き継ぎ実行される状態を示している。その際、構成制御処理部160は、共用データ記憶部10bに定義されている情報を参照してクラスタ管理機構900に対してアプリケーションの起動指示を出し(246)、指示を受け取ったクラスタ管理機構900は、232に示すように、アプリケーションAP1、AP2をそれぞれ起動する。構成制御処理部160は、226に示すように、クラスタ管理機構900から各サーバ及びアプリケーションの状態を取得する。

【0021】5. システム動作概要

次に、クラスタシステムの動作について説明する。図6は、クラスタシステムの動作概要図である。大きな流れは、以下のようになる。

(1) 常時実行の状況収集

常時、下記情報を取り込み、共用データ記憶部に書き込み、他サーバの共用データ記憶部に反映を行う。

(a) サーバ状態

状態監視部20は、クラスタにサーバが組み込まれているか否かの状態をクラスタ管理機構900より取り込み(矢印225)、全サーバの状態を管理する。サーバの突然の異常停止や異常停止後の立ち上がりなどを検知する。共用データ記憶部は、210に示すように、他サーバに反映を行う。

(b) 各サーバの負荷状態

各サーバの負荷更新部30a~30cは、自サーバの負荷をそれぞれ取り込み、共用データ記憶部へ書き込む。書き込まれた各サーバの負荷は、それぞれ他サーバへ反映される。他サーバへお互いに反映することにより、各サーバの共用データ記憶部は、全てのサーバの負荷状態を記憶する。

【0022】(2) アプリケーションのサーバ割り付けの管理、制御

構成制御処理部の総合的な管理として、構成制御部が下

記に示す各アプリケーションプログラムの実行状態の管理、制御を行う。

1) 管理

各サーバの状態と合わせて、各アプリケーションの下記状態を常に管理している。

(a) 各アプリケーションの動作サーバ

アプリケーションの動作サーバとは、そのアプリケーションが動作しているサーバである。

(b) 各アプリケーションのバックアップサーバ

各アプリケーションが動作しているサーバが停止した場合、どのサーバで再起動させるかを定義する情報である。再起動させるサーバをバックアップサーバという。

2) 制御

後述されるタイミング(構成サーバ変更、定周期)で、各アプリケーションの動作サーバとバックアップサーバを再定義し、クラスタ管理機構に動作サーバとバックアップサーバの再定義の指示を行う。なお、動作サーバ、バックアップサーバの再定義は、予め定義された制御動作定義ファイルに沿って行われる。制御動作定義ファイル35には、下記の2つの記述方式がある。

(a) 規約記述型

サーバ状態に応じた割り付け指示などの規約の詳細をシステム管理者が予め記述する。状況に応じた複雑な指示が可能という長所がある。だが、システム管理者の負荷が重いという欠点もある。

(b) モデル指定型

クラスタシステムが用意する代表的な形態のモデルをシステム管理者が指定する簡易な設定方式である。設定が簡単に短時間ででき、システム管理者の負荷を大幅に軽減するという長所がある。例えば、予め用意されたモデルから二重系ホットスタンバイ、多重系N対1バックアップ型などのモデルを指定できる。

【0023】(3) 初期立ち上げ時の動作

初期立ち上げ時には、構成制御部40が制御動作定義ファイルをもとに作成された共用データ記憶部の制御動作規約情報テーブル(後述)を参照して動作サーバとバックアップサーバを定義し、各アプリケーションの起動要求をクラスタ管理機構900に伝える。

【0024】(4) 構成サーバ変更時の動作

サーバの異常停止及び異常停止後の再立ち上げにより、構成サーバが変更した時、状態監視部20からのトリガ252で、各アプリケーションのバックアップサーバの再定義を構成制御部40が行う。構成制御部40は、再定義したバックアップサーバをクラスタ管理機構900に伝える。なお、バックアップサーバの再設定先は、サーバ構成及び各サーバの負荷状況をもとに決定される。

【0025】(5) 定周期による動作

負荷更新部30aからのトリガ254で、構成制御部40は、定期的に各サーバの負荷をチェックし、チェックの結果により、各アプリケーションのバックアップサーバ

バを低負荷のサーバに再設定し、アプリケーション移動時に特定サーバに負荷が集中することを防ぐ。

【0026】(6)サーバ異常時のアプリケーションの移動動作

実際に、サーバの異常停止が発生した場合、構成制御処理部160により設定されたバックアップサーバでアプリケーションを再起動させるために、クラスタ管理機構900にアプリケーションの再起動の指示を出す。指示を受け取ったクラスタ管理機構900がアプリケーションを再起動する。以上のように、この発明のクラスタシステムは動作する。

【0027】6.サーバ状態の監視

次に、状態監視部20が行うサーバ状態の監視について説明する。図7は、サーバ状態の監視を説明する図である。状態監視部20によるサーバ状態の取得は、このクラスタシステムが使用するクラスタ管理機構900の提供する状態取得機能を使用して行う。状態監視部20は、クラスタ管理機構900からポーリング形式でサーバ状態を取得する。本実施の形態では、今回取得したサーバ状態の情報と前回取得したサーバ状態の情報を比較し、状態変化を検知する方式を採用している。主たる動作は、以下の通りである。

(a)クラスタ管理機構900の提供するサーバ状態取得コマンドを使用し、ポーリング形式で全サーバの状態を取得する(矢印260)。

(b)初回の取得時、共用データ記憶部のサーバ状態情報テーブル12aに記憶し、他サーバのサーバ状態情報テーブル12b、12cへ反映を行う(矢印210)。

(c)2回目以降の取得時(前回の取得結果を記録済みの時)は、記録情報と比較する。状態変化を検知すると、共用データ記憶部のサーバ状態情報テーブル12aに取得した状態及び変更マーク(切り離し/組み込みマーク)を記憶し、他サーバのサーバ状態情報テーブル12b、12cへ反映を行う(矢印210)。

(d)初回及び変更時、構成制御部へ制御トリガを指示する。トリガは、OS(Operating System)が提供するメッセージ送信機能(プログラム間通信機能)を使用する。

【0028】7.負荷状況の監視

次に、負荷更新部30a~30cが行う負荷状況の監視について説明する。図8は、負荷状況の監視を説明する図である。各サーバの負荷更新部が自サーバの負荷を取り込み、共用データ記憶部10a~10cに書き込み、他サーバへ反映を行う。動作は、以下の通りである。

(a)負荷更新部30a~30cは、定周期で起動される。

(b)負荷更新部30a~30cは、OSが提供するCPU負荷情報取り込みコマンドで、OSツール群910a~910cより自サーバのCPU負荷を取り込み(262a~262c)、共用データ記憶部のサーバ負荷情

報テーブル13a~13cに書き込み(268a~268c)、他サーバの共用データ記憶部10a~10cへそれぞれ反映する。

(c)構成制御部40が動作するサーバ上で動作する負荷更新部の場合、構成制御部40へバックアップサーバチェックのトリガを送る(矢印264)。

【0029】8.アプリケーション状況の管理、制御次に、構成制御部40が行うアプリケーション状況の管理、制御について説明する。構成制御部40は、各アプリケーションの状態を管理し、状態監視部20又は負荷更新部30からのトリガによって、各アプリケーションのバックアップサーバの変更などを行う。

【0030】(1)管理データ

管理、制御するための必要な管理データは、以下のものである。なお、各データはすべて共用データ記憶部10a~10cに確保され、更新時は全てのサーバに反映される。従って、各サーバは同一イメージの管理データを有することになる。

【0031】(a)制御動作規約情報テーブル11

図9に、制御動作規約情報テーブル11を示す。制御動作規約情報テーブル11は、サーバの構成に応じ、各アプリケーションの動作サーバとバックアップサーバの割り付け定義が記述された動作規約定義である。定義自体は、定義ファイル制御動作定義ファイル35に記述され、定義ファイル制御動作定義ファイル35の内容が共用データ記憶部に制御動作規約情報テーブル11a~11cとしてクラスタの初期構成時に展開される。制御動作規約情報テーブル11は、クラスタシステム動作中に書き替えられることはなく、固定値である。各規約の意味は、以下の通りである。

・動作サーバ1105

クラスタの初期立ち上げ時にアプリケーションを動作させるサーバ名である。

・バックアップサーバ1107

初期立ち上げ時、設定されるバックアップサーバ名である。未定義時、バックアップサーバを割り付けない。

・固定指示1109

定期負荷によりバックアップサーバの切替を行うか否かの指定である。切り替えを行う時“OFF”、切り替えを行わない時“ON”とする。

・自動移動モード1111

同一サーバで複数のアプリケーションが動作している状態で、他のサーバが復旧した時、その復旧したサーバに複数のアプリケーションのうちのいずれかのアプリケーションを移動させるか否かの指定である。否の時、復旧したサーバへのアプリケーションの移動は行わずに、復旧したサーバへバックアップサーバを割り付ける処理となる。固定指示1109、自動移動モード1111の2つの規約は、アプリケーションを移動させるためには、そのアプリケーションを終了させなければならないの

で、業務の内容により途中で停止させることが望ましくない場合への対応として用意されている。

【0032】(b) サーバ状態情報テーブル12

図10に、サーバ状態情報テーブル12を示す。各サーバの正常／異常などの状態を管理するためのもので、状態監視部20が共用データ記憶部に作成、更新する。各パラメータの意味は、下記の通りである。

・サーバ状態1203

サーバがクラスタに組み込まれているか否か(正常／切り離し)の状態を示す。

・検出状態1205

状態変化のあったサーバについて、検出した状態の変化が“クラスタからの切り離し”か“クラスタへの組み込み”かを示す。

・処理モード1207

サーバ状態の変化に伴う変更処理が処理中か否かを示す。構成制御部160の構成制御動作途中にサーバの移動が発生した場合、サーバ移動後のリカバリに使用する。

【0033】(c) サーバ負荷情報テーブル13

図11に、サーバ負荷情報テーブル13を示す。各サーバのCPU負荷情報を管理するためのもので、負荷更新部が定期的に更新する。サーバ負荷情報テーブル13は、共用データ記憶部に確保される。各パラメータの意味は、下記の通りである。

・負荷情報1303

各サーバのCPU負荷(%)を示す。

【0034】(d) アプリケーション動作情報テーブル14

図12に、アプリケーション(AP)動作情報テーブル14を示す。各アプリケーションの動作サーバとバックアップサーバを管理するもので、構成制御部40が作成、更新する。また、構成制御部40は、アプリケーション動作情報テーブル14の更新内容に従って、クラスタ管理機構900に実行指示を行う。各パラメータの意味は、下記の通りである。

・動作サーバ名1403

アプリケーションが実際に動作しているサーバ名が保持される。

・バックアップサーバ名1405

実際にバックアップサーバとして割り付けられているサーバ名が保持される。

・処理中フラグ1407

アプリケーションに対する動作サーバやバックアップサーバの変更処理が、処理中か否かを示す。構成制御部160の構成制御動作途中に、サーバの移動が発生した場合、サーバ移動後のリカバリに使用する。

【0035】(2) 構成制御部40の動作トリガ

構成制御部40の動作トリガには、下記のものがある。

・初期立ち上げ

システムが開始されたタイミングで、状態監視部20から通知される。

・構成変更時

サーバ異常停止、再立ち上げ、操作指示によるサーバ切り離し／組み込みにより、サーバの状態が変更されたタイミングで、状態監視部20から通知される。

・定期

負荷更新部が定期的に行う自サーバ負荷更新タイミングで、構成制御部40が動作しているサーバの負荷更新部30から通知を受けとる。構成制御部40は、動作トリガの通知をOS提供のメッセージ通知機構(プログラム間通信機能)で受けとる。構成制御部40は、メッセージの通知を常に待ち、受けとると要因に沿った処理を行う。

【0036】(3) 初期立ち上げ時の構成制御部40の動作

図13は、構成制御部40の初期立ち上げ時の動作を示す流れ図である。構成制御部40は、状態監視部20から何らかのメッセージの送付を受けると、S101において、メッセージを取り込み、メッセージの要因が“初期立ち上げ”かどうか判定する。NOの場合、S103において、他の要求のメッセージかどうかチェックを行う。YESの場合、S105において、サーバ状態情報テーブルから現状のサーバ構成を取得し、本構成に対応する制御動作規約情報テーブル11を参照し、各アプリケーションの動作サーバとバックアップサーバを決定し、アプリケーション動作情報テーブルを作成する。次に、S107で、クラスタ管理機構900に、各アプリケーションをアプリケーション動作情報テーブルに定義された動作サーバで起動するよう指示を行う。次に、S109において、クラスタ管理機構900にアプリケーションのバックアップサーバの設定を行うよう指示を出す。S111において、全アプリケーションの処理が終了したか否かを判定し、終了するまでS105～S109の処理を繰り返す。このように、全アプリケーションに対し、上記起動／設定処理を行う。S111の判定で全アプリケーション処理が終了すると、S113でメッセージ待ちに戻る。この状態で、アプリケーション動作サーバが異常停止すると、クラスタ管理機構900は、構成制御部40の指示により設定されているバックアップサーバでアプリケーションを再起動する。

【0037】(4) 構成変更時の構成制御部40の動作

図14から図18は、サーバ構成が変更された時の構成制御部40の動作を示す流れ図である。

1) 動作サーバ切り離し時

サーバの異常停止やユーザ指示によりサーバがクラスタから削除された場合、構成制御部40は、図14、図15に示した流れ図に従い、下記動作を行う。構成制御部40は、状態監視部20からメッセージの送付を受けると、S121において、送付されたメッセージを取り込

み、取り込んだメッセージが“切り離し要求”かどうか判定する。判定した結果、“切り離し要求”でなければ、S123において、他の要因のチェックを行う。“切り離し要求”であった場合には、S125において、サーバ状態情報テーブル12から切り離し要求のあるサーバ(検出状態1205が切り離しのサーバ)をサーチする。検出したら、サーバ状態情報テーブル12の処理モード1207を処理中にセットする(S129)。検出しない場合は、元のメッセージ待ちに戻る(S127)。

【0038】次に、S131において、アプリケーション動作情報テーブルから該当サーバで動作していたアプリケーションをスキャンする。次に、S133において、検出したアプリケーションのアプリケーション動作情報テーブルの処理中フラグ1407をセットし、実際に動作しているサーバ名をクラスタ管理機構900から取り込み、取り込んだサーバ名をアプリケーション動作情報テーブル14の動作サーバ名1403にセットする。この時点では、切り離し要求のあるサーバ上で動作していたアプリケーションは、既にクラスタ管理機構900によりアプリケーション動作情報テーブル14で定義されているバックアップサーバで再起動しているはずであるが、構成制御部40は、クラスタ管理機構900から実際の動作サーバ名を取り込む。

【0039】実際に動作しているサーバが存在しない(アプリケーション停止=バックアップ未定義又は重複障害)場合は、図15のS139に示すように、バックアップサーバ無しを選択する。

【0040】アプリケーションが動作中の場合、制御動作規約情報テーブルの同一サーバ構成の項を参照し(図14、S135)、アプリケーションのバックアップサーバが未定義の場合、図15のS139に示すように、バックアップサーバ無しを選択する。アプリケーションのバックアップサーバの未定義は、故障により動作サーバ数が特定台数に満たない場合、サーバの処理能力上の問題から、重要でないアプリケーションを切り捨て、重要なアプリケーションの動作を保証するため等に使用される。

【0041】バックアップサーバが定義されている場合、アプリケーション動作情報テーブル/サーバ状態情報テーブルから、アプリケーションが動作していない正常サーバをサーチする(図15、S137)。検出した場合、検出した正常サーバをバックアップサーバに決定する(図15、S143)。

【0042】検出しない場合、サーバ負荷テーブルから最も低負荷なサーバをバックアップサーバに選択する(図15、S141)。

【0043】次に、S145で、アプリケーション動作情報テーブルのバックアップサーバ名を選択したバックアップサーバに変更し、他サーバに反映を行う。その

後、S147で、クラスタ管理機構900にバックアップサーバの変更指示を行う。完了後、アプリケーション動作情報テーブル14の処理中フラグ1407をリセットし、他サーバへ反映を行う(S149)。S131～S149の処理を該当する全アプリケーションに対して繰り返し行う。

【0044】全アプリケーションに対する処理が完了したら、サーバ状態情報テーブル12の該当サーバのサーバ状態1203を切り離し状態にし、検出状態1205をクリア、処理モード1207をリセットし、他サーバへ反映する(S151)。

【0045】2)サーバ再組み込み時

異常サーバの再立ち上げやコマンド指示によるサーバ組み込みで、クラスタにサーバが再組み込みされた場合、構成制御部40は、図16から図18に示す流れ図に従い、下記動作を行う。

【0046】構成制御部40は、状態監視部20からメッセージを送付されると、S161において、送付されたメッセージを取り込み、メッセージが“組み込み要求”かどうか判定する。“組み込み要求”でなかった時には、S163において、他の要求のチェックを行う。“組み込み要求”であった場合には、S165において、サーバ状態情報テーブル12から組み込み要求のあるサーバ(検出状態1205が組み込み)をサーチする。検出したら、サーバ状態情報テーブル12の処理モード1207を処理中にセットする(S169)。検出しない場合は、元のメッセージ待ちに戻る(S167)。

【0047】次に、アプリケーション動作情報テーブル14をサーチし、動作していないアプリケーション(動作サーバ名1403に登録のないアプリケーション)をサーチする(S171)。検出した場合、下記を行う。アプリケーション動作情報テーブル14の処理中フラグ1407を処理中にし、動作サーバ名1403を組み込みサーバに設定し、他サーバへ反映する(S173)。次に、S175で、アプリケーションの起動指示をクラスタ管理機構900に対して行う。完了後、アプリケーション動作情報テーブル14の処理中フラグ1407をリセットし、他サーバへ反映する。該当するアプリケーションに対し処理を行ったら、元のメッセージ待ちへ戻る(S179)。

【0048】S171の判定で、動作サーバなしのアプリケーションが無い場合は、図17のS181において、同一サーバで重複動作しているアプリケーションをアプリケーション動作情報テーブル14からサーチする。存在する場合、下記のアプリケーション移動処理を行う。制御動作規約情報テーブル11を参照し、同一サーバで重複動作しているアプリケーションの内、同一サーバ構成時、自動移動モードがONのものをサーチする(S183)。存在しない場合は、アプリケーション移

動処理はスキップし、図18のS193に進む。検出した場合、図17のS185において、そのアプリケーションのアプリケーション動作情報テーブル14の処理中フラグ1407をセットし、動作サーバ名1403を組み込みサーバに変更し、他サーバに反映する。その後、クラスタ管理機構900に、該当アプリケーションの移動(停止、起動サーバ設定、起動)指示を行う(S187)。完了後、アプリケーション動作情報テーブル14の処理中フラグ1407をリセットし、他サーバに反映する(S189)。処理終了後、元のメッセージ待ちへ戻る(S191)。

【0049】上記のアプリケーション移動処理が行われなかった場合、図18のS193において、各アプリケーションのバックアップサーバ名1405を組み込みサーバに変更する。本処理は、前述同様、アプリケーション動作情報テーブル14の処理中フラグ1407を処理中にし、バックアップサーバ名1405を書き換え、他サーバに反映する。次に、S195において、クラスタ管理機構900に、バックアップサーバ名変更指示を行う。次に、S197において、アプリケーション動作情報テーブル14の処理中フラグ1407をリセットし、他サーバに反映する。

【0050】最後に、サーバ状態情報テーブル12のサーバ状態1203を正常に書き換え、検出状態1205をクリア、処理モード1207をリセットし、他サーバへ反映する(S199)。

【0051】(5) 定期負荷更新時の動作
負荷更新部からのトリガがあった場合、構成制御部40は、図19に示す流れ図に従い、下記動作を行う。構成制御部40は、負荷更新部30からメッセージを送付されると、S211で送付されたメッセージを取り込み、
“定期負荷要求”かを判定する。“定期負荷要求”でない時、S213で他の要求のチェックを行う。“定期負荷要求”であった時、S215でアプリケーション動作情報テーブル14をサーチし、バックアップサーバが定義されている動作アプリケーションの内、制御動作規約情報テーブルの同一サーバ構成時の固定バックアップ指定なしのアプリケーション(固定指示1109がOFFのアプリケーション)をサーチする。

【0052】該当するアプリケーションがサーチされると、S217でサーチされたアプリケーションのバックアップサーバの負荷と他サーバの負荷をサーバ負荷情報テーブルで比較しながら、より低負荷のサーバをサーチする。また、サーチしたより低負荷のサーバをサーバ状態情報テーブル12でサーチし、サーバ状態1203が正常であることを確認する。より低負荷で正常なサーバがあれば、S219で、該当アプリケーションのアプリケーション動作情報テーブル14の処理中フラグ1407を処理中にし、バックアップサーバ名1405をサーチしたサーバに書き換え、他サーバに反映する。次に、

S221でクラスタ管理機構900にバックアップサーバ名の変更指示を行う。完了後、アプリケーション動作情報テーブル14の処理中フラグ1407をリセットし、他サーバに反映する。(S223)。全アプリケーションに対し、上記S215～S223の処理を行い(S225)、処理が終了したらメッセージ待ちへ戻る(S227)。

【0053】9. 制御動作例

以下に、制御動作例として、サーバ状態の変更に対するアプリケーションのバックアップサーバの割り付けの例をいくつか記述する。

(1) 構成変更(異常停止/復旧)による動作例

5 台構成時の構成変更時の割り付け例を示す。図20は、サーバ構成の変更例を示す図である。図21は、図20のサーバ構成の変更にそれぞれ対応するアプリケーションの割り付け例を示す図である。図20の2001、図21の2101に示すように、各アプリケーションAP1、AP2、AP3がサーバ100a、100b、100cで動作し、各アプリケーションのバックアップサーバは、サーバ100dに割り当てられている。

【0054】次に、図20の2002に示すように、サーバ100aが異常停止すると、クラスタ管理機構900がアプリケーションAP1を、2101においてバックアップサーバに指定されているサーバ100dに移す。構成制御部40は、このトリガを元に、各アプリケーションのバックアップサーバをアプリケーションが何も動作していないサーバ100eに割り付ける。割り付けた状態を図21の2102に示す。

【0055】次に、図20の2003に示すように、サーバ100bが異常停止すると、クラスタ管理機構900がアプリケーションAP2をバックアップサーバに指定されているサーバ100eに移す。構成制御部40は、このトリガを元に、各アプリケーションのバックアップサーバをサーバ100c、100d、100dと、図21の2103に示すように再設定する。負荷による調整は、前述した説明のように行われるので、ここでは説明は省略する。

【0056】次に、図20の2004に示すように、サーバ100aが復旧して立ち上がり、それを構成制御部40が検知すると、このトリガを元に、構成制御部40は、各アプリケーションのバックアップサーバを、図21の2104に示すように、サーバ100aに再設定する。

【0057】(2) 構成変更(復旧)による動作例

図22から図24は、異常縮退後の復旧時の動作例を示す図である。図22の2201、図23の2301に示すように、サーバ100d、100eが正常な状態で稼働しており、アプリケーションAP1、アプリケーションAP3がサーバ100dで重複動作しているものとする。その状態で、サーバ100cが復旧すると、構成制

御部40は、そのトリガを受け、サーバ100c, 100d, 100eの3台によるクラスタ構成時の制御動作規約情報テーブルのアプリケーションAP3の自動移動モードを判定する。自動移動モードONの場合、構成制御部40は、クラスタ管理機構900に指示を出し、アプリケーションAP3はサーバ100cに移動(停止、再起動)される。その結果、図22の2202、図23の2303に示すように、サーバ100c, 100d, 100eでの運転に移行する。

【0058】なお、自動移動モードOFFの場合は、アプリケーションの移動は行われず、各アプリケーションのバックアップサーバにサーバ100cが割り当てられる処理となる。この時の割り当てを、図24に示す。

【0059】(3) CPU負荷による動作例

図25から図28は、CPU負荷変動による定期更新の動作例を示す図である。図25に示すように、サーバ100c, 100d, 100eでアプリケーションAP3, AP1, AP2がそれぞれ動作している。各CPU負荷により図26に示すように、バックアップサーバがサーバ100c, 100d, 100dと割り付けられている。サーバ100dで動作しているアプリケーションAP1のバックアップは、サーバ100eより軽負荷であるサーバ100cに割り付けられている。負荷更新部30による負荷の更新後、図27に示すように、サーバ負荷の軽負荷順位で、サーバ100cとサーバ100eが逆転したので、構成制御部40は、サーバ100cをバックアップサーバとしていたアプリケーションAP1のバックアップサーバを図28に示すように、サーバ100eに切り替える。

【0060】10. 制御動作規約の定義

次に、制御動作規約の定義について、図2.9を用いて説明する。共用データ記憶部10の制御動作規約情報テーブル11作成のために、ユーザが記述する定義ファイル2901について説明する。制御動作規約情報テーブル11の作成方式は、規約記述型とモデル指定型の2方式がある。

【0061】a) 規約記述型

サーバ構成の各変動に対し、どのような動作制御を行うかの規約定義2902を詳細に記述したユーザ定義ファイル2901を図29の2903に示すように、ユーザが記述して作成する方式である。このユーザ定義ファイル2901は、プログラミングを必要としない簡易なビルドインコマンドで記述可能とする。クラスタシステム立ち上げ時に構成制御部40により規約定義2902が参照され、制御動作規約情報テーブル11が作成される。

【0062】b) モデル指定型

特定の動作モデルを予め定義し、その中からユーザに指定をさせる方式である。当然、内部で、指定されたモデルをもとにして、規約定義2916に変換が行われる。

記述できるモデルは、以下の通りである。

・二重系ホットスタンバイ(ロードシェア含む)型モデル

図30に、二重系ホットスタンバイモデルの定義例を示す。このモデルでは、アプリケーションを双方のサーバで動くよう2個用意した場合、ロードシェア型となる。アプリケーション動作中のサーバが異常停止した場合、動作中であったアプリケーションは、もう一方のサーバで再起動される。異常サーバが復旧した場合、復旧したサーバは、各アプリケーションのバックアップサーバとなる。

・多重系N対1バックアップ(バックアップ浮動)型

図31に、多重系N対1バックアップモデルの定義例を示す。(N+1)台のサーバで構成され、1台を共通のバックアップとする形態である。アプリケーションはN個存在し、N台のサーバでそれぞれ動作する。動作サーバが停止すると、アプリケーションはバックアップサーバで動作する。異常サーバが復旧し、再度、組み込まれると、全アプリケーションのバックアップサーバとなる。なお、2個以上のサーバが異常停止した場合、そのアプリケーションは切り捨てられ、サーバが復旧した時、再起動される。どちらの作成方式の場合も、記述ファイル内容はクラスタシステム立ち上げ時、図29の2905, 2918に示すように、制御動作規約情報テーブル11に展開される。

【0063】以上のように、この実施の形態では、クラスタ(分散した複数の計算機を一まとまりにしたシステム)構成に於いて、計算機及びアプリケーションの状態変化に応じて、アプリケーションの起動、停止、起動計算機の決定、及び引継ぎ計算機(アプリケーション、もしくはアプリケーション実行中の計算機が異常停止した場合、次にアプリケーションを実行する計算機のこと)の決定を自動的、かつ、最適に制御できるようにしたクラスタの自動構成制御方式について説明した。この実施の形態のクラスタシステムによれば、計算機及びアプリケーションの状態を常に監視し、その状態の変化に応じて、自動的かつ最適にアプリケーションの動作計算機を動的に制御するので、システム管理者の負担が軽減される。

【0064】前述した実施の形態では、ネットワークを介して更新される共用データ記憶部をサーバ毎に備える場合について説明したが、クラスタシステム内で共用のメモリ装置を備えて共用データ記憶部としてもよい。

【0065】

【発明の効果】この発明によれば、クラスタシステムを構成するサーバの状態に応じて、アプリケーションとサーバとの対応を動的に制御できる。これにより、システム管理者がシステム監視、制御に要する負担を大幅に軽減し、効率のよい高信頼システムの構築が可能になる。

【0066】また、この発明によれば、クラスタに構成

されている計算機のメンバーに応じたアプリケーションの引継ぎ計算機を自動的に決定できる。

【0067】また、この発明によれば、新たにサーバが起動した時、操作者を介入させずアプリケーションを自動引継ぎできる。

【0068】また、クラスタのメンバーに変動が生じた際、実行中のアプリケーションの引継ぎ計算機を変更し、決定するので、常に最適な引継ぎ計算機を選択可能である。

【0069】また、サーバの負荷に応じて、アプリケーションの引き継ぎを行うので、負荷の分散が図れる。

【0070】また、定義ファイルを基にクラスタの構成を制御する規約情報を作成するので、プログラミングレスによる定義が可能である。

【0071】また、システムで定義テンプレートを用意するので、パターン化された処理の場合、ユーザは、クラスタシステムの詳細な定義を行う必要がない。

【0072】また、構成制御部が、定期的に取得された各計算機の負荷(CPU、メモリ)状況によって最適な引継ぎ計算機を自動決定し、引継ぎを行うので、操作者が介入せず最適な引継ぎができる。

【図面の簡単な説明】

【図1】 この発明のクラスタシステムの構成を示すブロック図である。

【図2】 この発明のクラスタシステムの構成を示すブロック図である。

【図3】 この発明のクラスタシステムのハードウェア構成図である。

【図4】 この発明のクラスタシステムの機能を示すブロック図である。

【図5】 この発明のクラスタシステムの機能を示すブロック図である。

【図6】 この発明のクラスタシステムの動作概要を示す図である。

【図7】 この発明のクラスタシステムの状態監視部20が行うサーバ状態の監視を説明する図である。

【図8】 この発明のクラスタシステムの負荷更新部30が行う負荷状況の監視を説明する図である。

【図9】 この発明のクラスタシステムの共用データ記憶部10の制御動作規約情報テーブル11を示す図である。

【図10】 この発明のクラスタシステムの共用データ記憶部10のサーバ状態情報テーブル12を示す図である。

【図11】 この発明のクラスタシステムの共用データ記憶部10のサーバ負荷情報テーブル13を示す図である。

【図12】 この発明のクラスタシステムの共用データ記憶部10のアプリケーション動作情報テーブル14を示す図である。

【図13】 この発明のクラスタシステムの構成制御部40の初期立ち上げ時の動作の流れ図である。

【図14】 この発明のクラスタシステムの構成制御部40のサーバ構成変更時の動作を示す流れ図である。

【図15】 この発明のクラスタシステムの構成制御部40のサーバ構成変更時の動作を示す流れ図である。

【図16】 この発明のクラスタシステムの構成制御部40のサーバ構成変更時の動作を示す流れ図である。

【図17】 この発明のクラスタシステムの構成制御部40のサーバ構成変更時の動作を示す流れ図である。

【図18】 この発明のクラスタシステムの構成制御部40のサーバ構成変更時の動作を示す流れ図である。

【図19】 この発明のクラスタシステムの構成制御部40の負荷更新時の動作の流れ図である。

【図20】 この発明のクラスタシステムのサーバ構成の変更例を示す図である。

【図21】 この発明のクラスタシステムの図20のサーバ構成の変更に対応するアプリケーションの割り当て例を示す図である。

【図22】 この発明のクラスタシステムの異常縮退後の復旧時の構成の変更例を示す図である。

【図23】 この発明のクラスタシステムの図22に対応するサーバの異常縮退後の復旧時のアプリケーションの割り当て例を示す図である。

【図24】 この発明のクラスタシステムの図22に対応するサーバの異常縮退後の復旧時のアプリケーションの割り当て例を示す図である。

【図25】 この発明のクラスタシステムのCPU負荷変動による定期更新時の動作例を示す図である。

【図26】 この発明のクラスタシステムのCPU負荷変動による定期更新時の動作例を示す図である。

【図27】 この発明のクラスタシステムのCPU負荷変動による定期更新時の動作例を示す図である。

【図28】 この発明のクラスタシステムのCPU負荷変動による定期更新時の動作例を示す図である。

【図29】 この発明のクラスタシステムの共用データ記憶部10の制御動作規約情報テーブル11の作成を説明する図である。

【図30】 この発明のクラスタシステムの二重系ホットスタンバイモデルの定義例を示す図である。

【図31】 この発明のクラスタシステムの多重系N対1バックアップモデルの定義例を示す図である。

【符号の説明】

10, 10a, 10b, 10c 共用データ記憶部、11, 11a, 11b, 11c 制御動作規約情報テーブル、12, 12a, 12b, 12c サーバ状態情報テーブル、13, 13a, 13b, 13c サーバ負荷情報テーブル、14, 14a, 14b, 14c アプリケーション動作情報テーブル、20 状態監視部、30, 30a, 30b, 30c 負荷更新部、35 制御動作

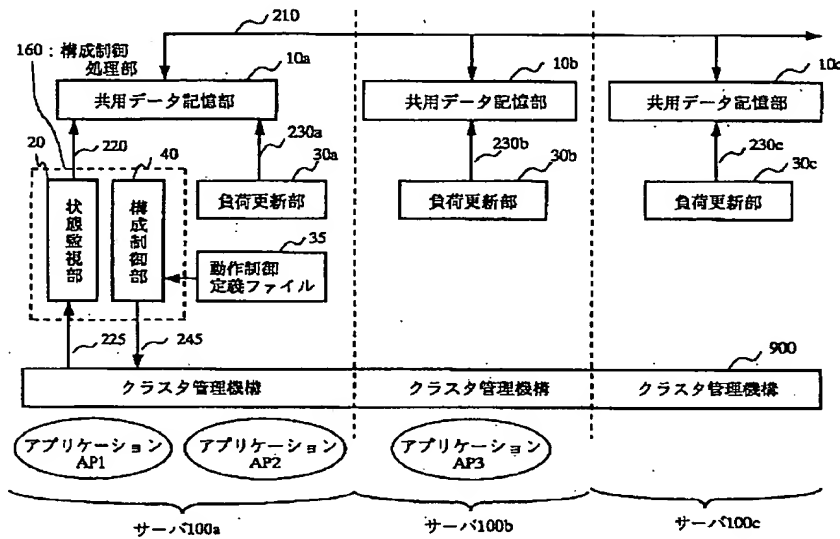
21

定義ファイル、40 構成制御部、100a, 100b, 100c, 100n サーバ、160 構成制御処理部、200a, 200b, 200n ディスク、900 クラスタ管理機構、910a, 910b, 910c OSツール群、1000 ネットワーク、1101 サーバ構成、1103 アプリケーション名、1105 動作サーバ、1107 バックアップサーバ、1109

22

固定指示、1111 自動移動モード、1201 サーバ名、1203 サーバ状態、1205 検出状態、1207 処理モード、1301 サーバ名、1303 負荷情報、1401 アプリケーション名、1403 動作サーバ名、1405 バックアップサーバ名、1407 処理中フラグ、AP1, AP2, AP3 アプリケーション。

【 図1 】

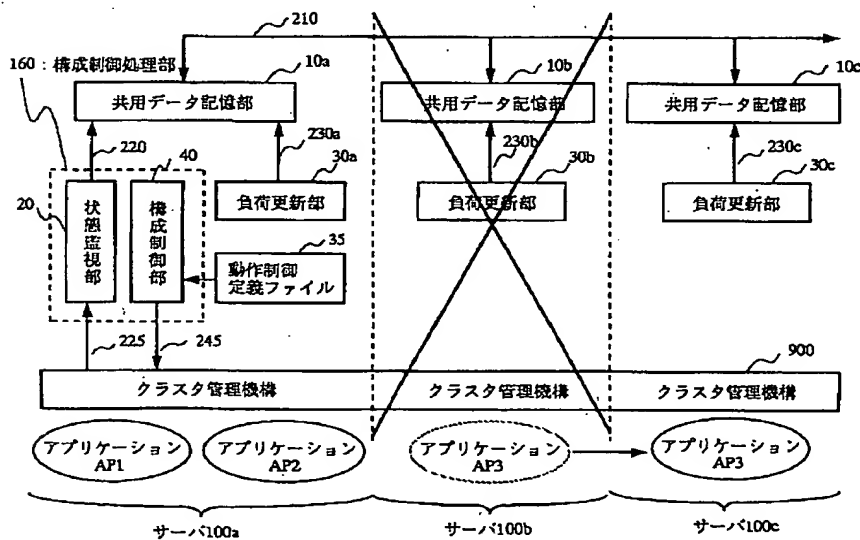


【 図11 】

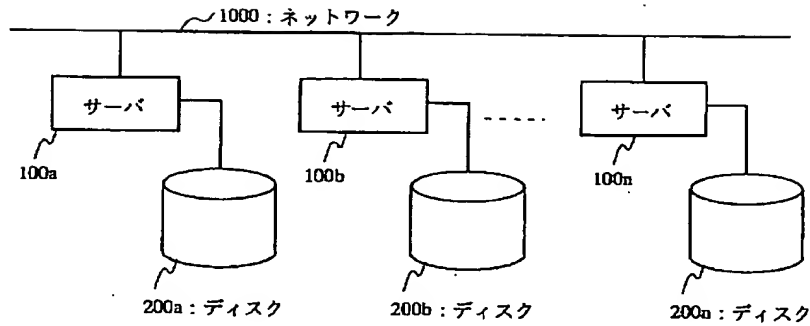
サーバ負荷情報テーブル 13

サーバ名	負荷情報
サーバ100a	0%
サーバ100b	40%
サーバ100c	60%
サーバ100d	20%

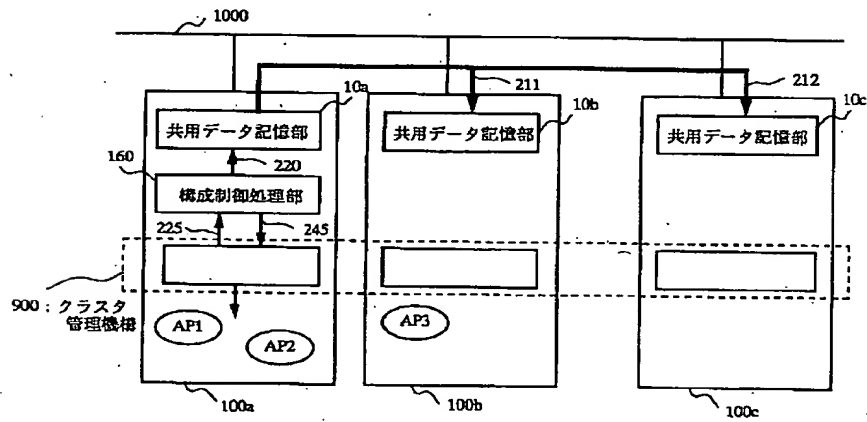
【 図2 】



【 図3 】



【 図4 】



【 図5 】

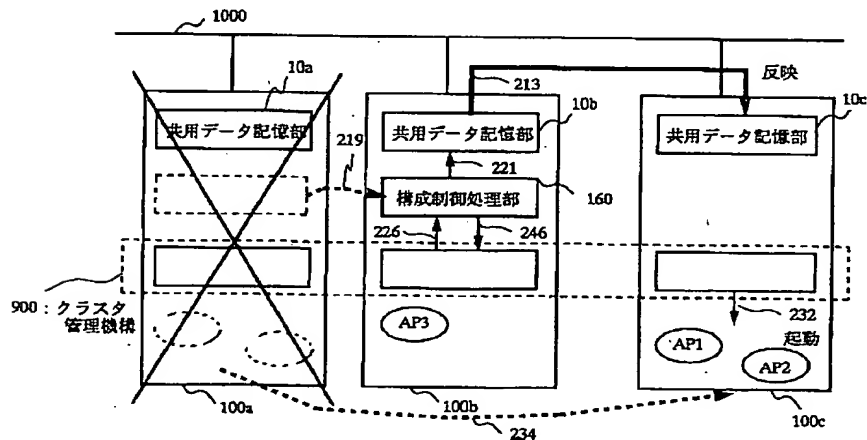


Figure 1 is a configuration diagram of a cluster management system. It shows a cluster management device (900) at the bottom, which is connected to a status display section (20) via a bus (260). The status display section (20) sends update requests (更新) to a shared data storage section (10a). The shared data storage section (10a) is connected to three server status information tables (12a, 12b, 12c) via a common data bus (210). The shared data storage section (10a) is also connected to a common data storage section (10b) and a common data storage section (10c).

[illegible]

【 図9 】

制御動作規約情報テーブル 11

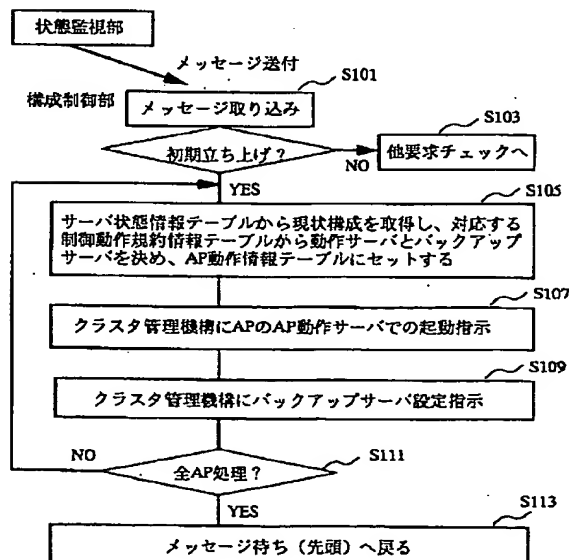
サーバ構成	AP名 1101	動作サーバ 1103	バックアップサーバ 1105	固定指示 1107	自動移動モード 1109
サーバ100a、サーバ100b サーバ100c、サーバ100d	AP1	サーバ100a	サーバ100d	ON	OFF
	AP2	サーバ100b	サーバ100d	ON	OFF
	AP3	サーバ100c	サーバ100d	ON	OFF
サーバ100b、 サーバ100c、サーバ100d	AP1	サーバ100d	サーバ100c	OFF	ON
	AP2	サーバ100b	サーバ100d	OFF	ON
	AP3	サーバ100c	サーバ100d	OFF	ON
サーバ100c、サーバ100d	AP1	サーバ100d	サーバ100c	OFF	ON
	AP2	サーバ100d	サーバ100c	OFF	ON
	AP3	サーバ100c	サーバ100d	OFF	ON
			⋮		

【 図10 】

サーバ状態情報テーブル 12

サーバ名 1201	サーバ状態 1203	検出状態 1205	処理モード 1207
サーバ100a	切り離し		
サーバ100b	正常	異常停止	処理中
サーバ100c	正常		
サーバ100d	正常		

【 図13 】

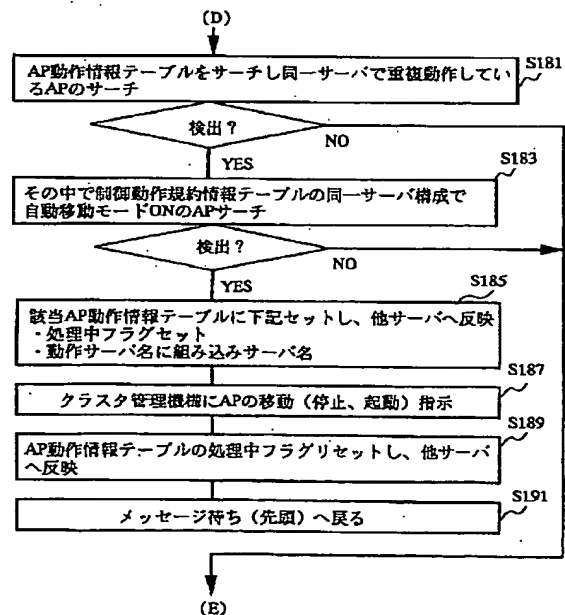


【 図12 】

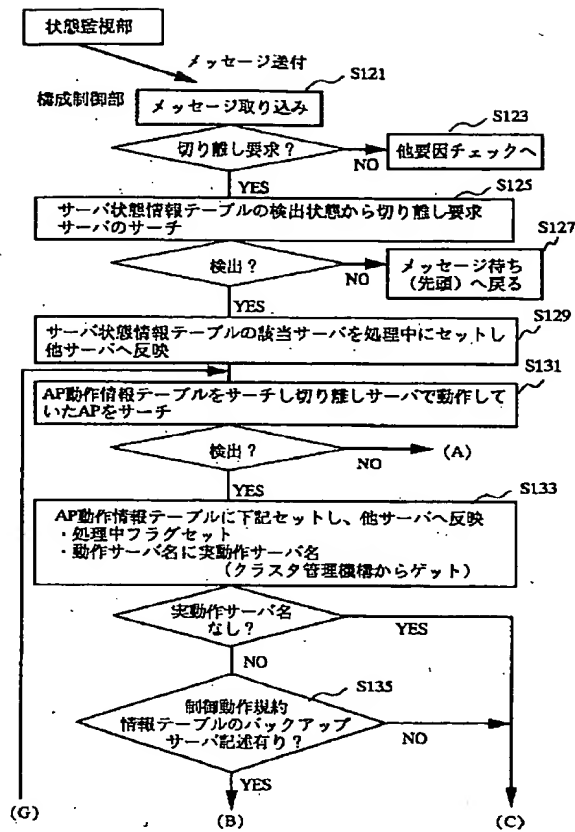
AP動作情報テーブル 14

AP名 1401	動作サーバ名 1403	バックアップサーバ名 1405	処理中フラグ 1407
AP1	サーバ100d	サーバ100c	完了
AP2	サーバ100b	サーバ100d	処理中
AP3	サーバ100c	サーバ100d	

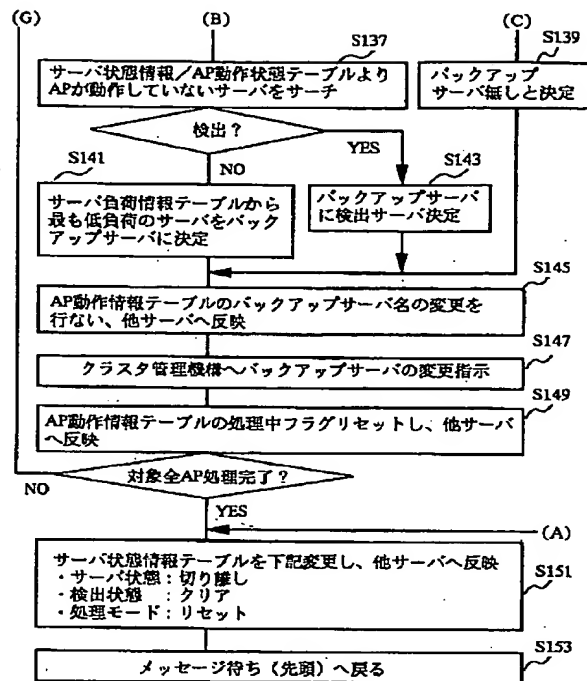
【 図17 】



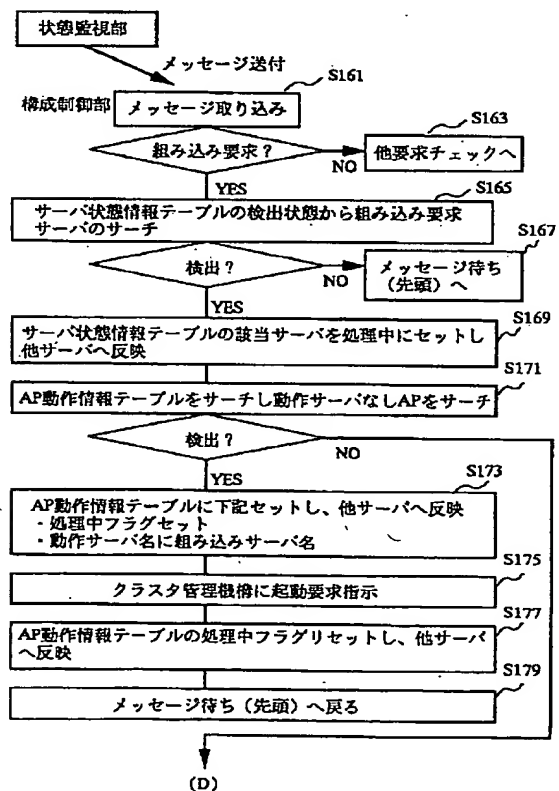
【 図14 】



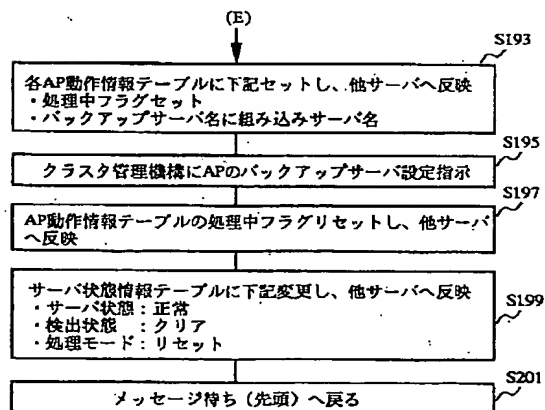
【 図15 】



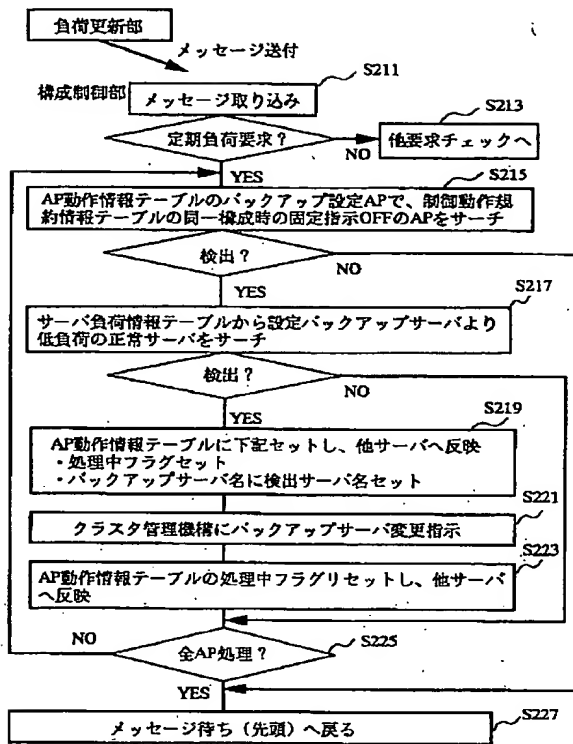
【 図16 】



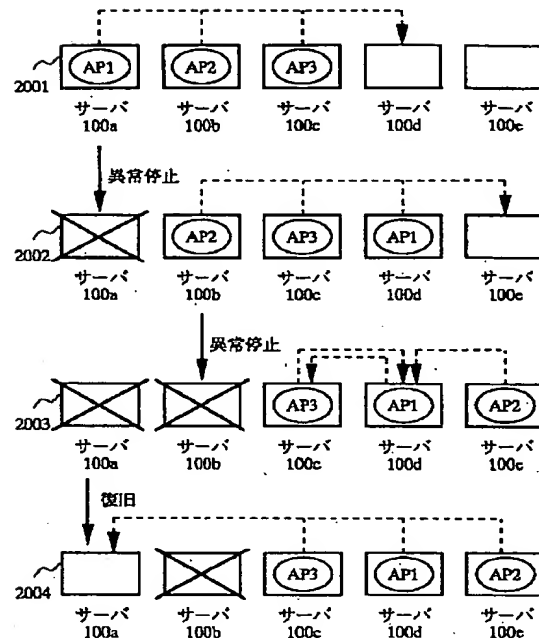
【 図18 】



【 図19 】



【 図20 】



【 図21 】

AP名	動作サーバ	バックアップサーバ
AP1	サーバ100a	サーバ100d
AP2	サーバ100b	サーバ100d
AP3	サーバ100c	サーバ100d

2101

AP名	動作サーバ	バックアップサーバ
AP1	サーバ100d	サーバ100c
AP2	サーバ100b	サーバ100c
AP3	サーバ100c	サーバ100c

2102

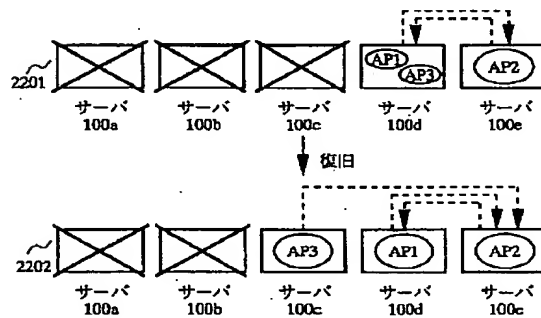
AP名	動作サーバ	バックアップサーバ
AP1	サーバ100d	サーバ100c
AP2	サーバ100c	サーバ100d
AP3	サーバ100c	サーバ100d

2103

AP名	動作サーバ	バックアップサーバ
AP1	サーバ100d	サーバ100a
AP2	サーバ100c	サーバ100a
AP3	サーバ100c	サーバ100a

2104

【 図22 】



【 図24 】

AP名	動作サーバ	バックアップサーバ
AP1	サーバ100d	サーバ100c
AP2	サーバ100c	サーバ100c
AP3	サーバ100d	サーバ100c

【 図 2 3 】

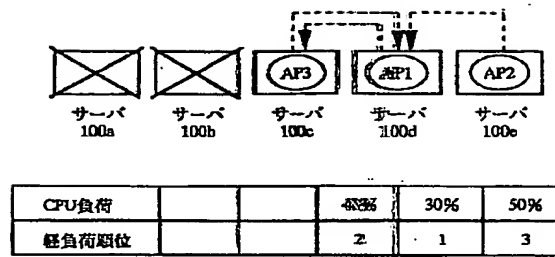
AP名	動作サーバ	バックアップサーバ
AP1	サーバ100d	サーバ100e
AP2	サーバ100e	サーバ100d
AP3	サーバ100d	サーバ100e

2301

AP名	動作サーバ	バックアップサーバ
AP1	サーバ100d	サーバ100e
AP2	サーバ100e	サーバ100d
AP3	サーバ100e	サーバ100e

2303

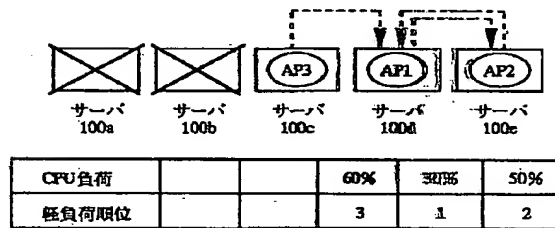
【 図 2 5 】



【 図 2 6 】

AP名	動作サーバ	バックアップサーバ
AP1	サーバ100d	サーバ100e
AP2	サーバ100e	サーバ100d
AP3	サーバ100e	サーバ100d

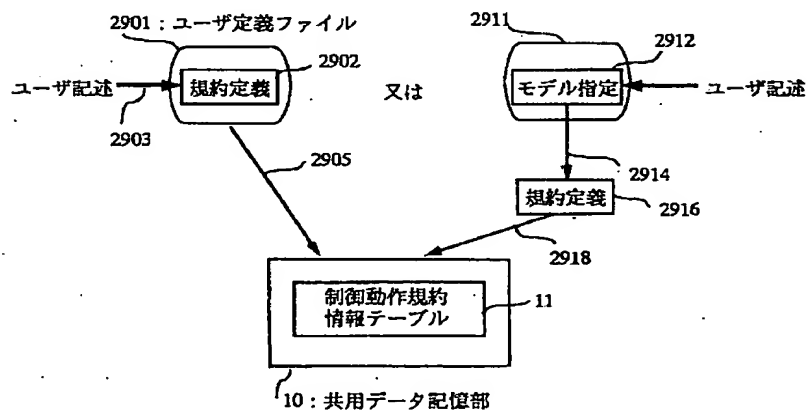
【 図 2 7 】



【 図 2 8 】

AP名	動作サーバ	バックアップサーバ
AP1	サーバ100d	サーバ100e
AP2	サーバ100e	サーバ100d
AP3	サーバ100e	サーバ100d

【 図 2 9 】



【 図30 】

モデル指定型：二重系ホットスタンバイモデルの定義ファイル

サーバ構成	AP名	動作サーバ	バックアップ サーバ	固定指示	自動移動 モード
サーバ100a、 サーバ100b	AP1	サーバ100a	サーバ100b	OFF	ON
	AP2	サーバ100b	サーバ100a	OFF	ON

【 図31 】

モデル指定型：多重系N対1定義ファイル

サーバ構成	AP名	動作サーバ	バックアップ サーバ	固定指示	自動移動 モード
サーバ100a、 サーバ100b、 サーバ100c	AP1	サーバ100a	サーバ100c	OFF	ON
	AP2	サーバ100b	サーバ100c	OFF	ON